

Introducción a la pan-genómica microbiana

Introducción a la filoinformática – pan-genómica y filogenómica microbiana,
TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>

Microbial genome evolution - the pan-genome

Pablo Vinuesa
Centro de Ciencias Genómicas – UNAM
vinuesa @ ccg_unam_mx
<http://www.ccg.unam.mx/~vinuesa/>

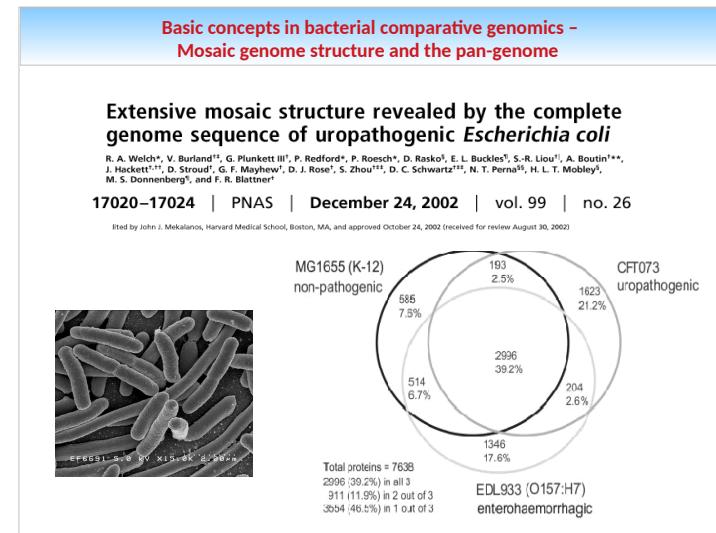
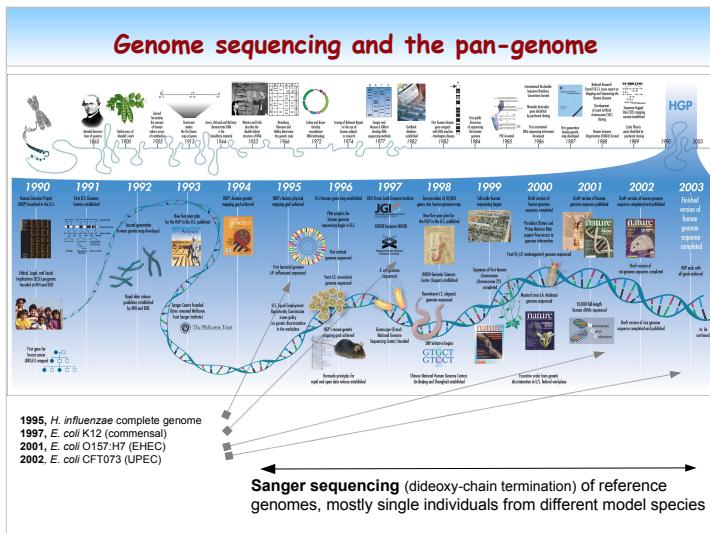
 @pvinmex

Introducción a la Filoinformática:
Pan-genómica y Filogenómica microbiana –
NNB & CCG-UNAM, 1-5 Agosto 2022
<https://github.com/vinuesa/TIB-filoinfo>

Genome evolution and the pan-genome

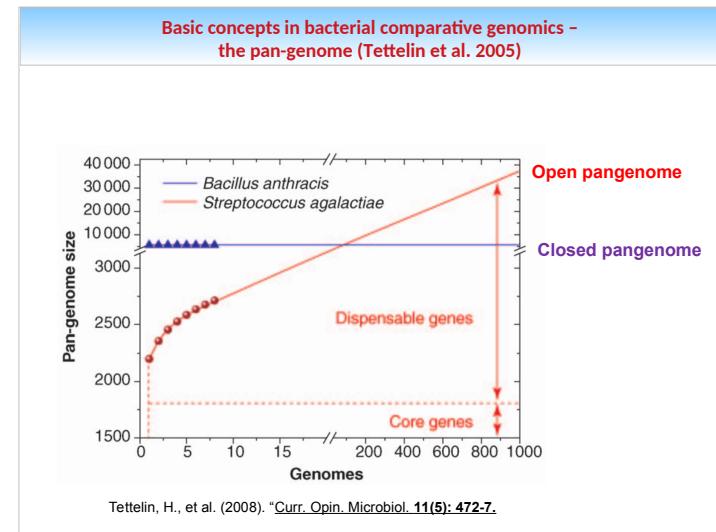
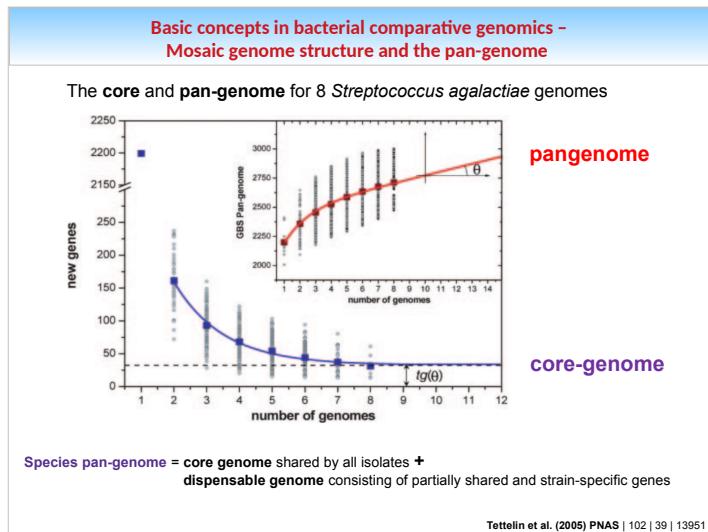
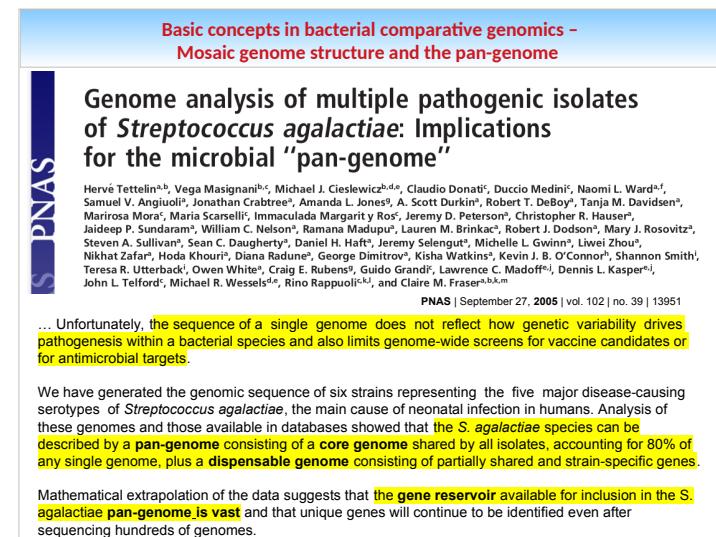
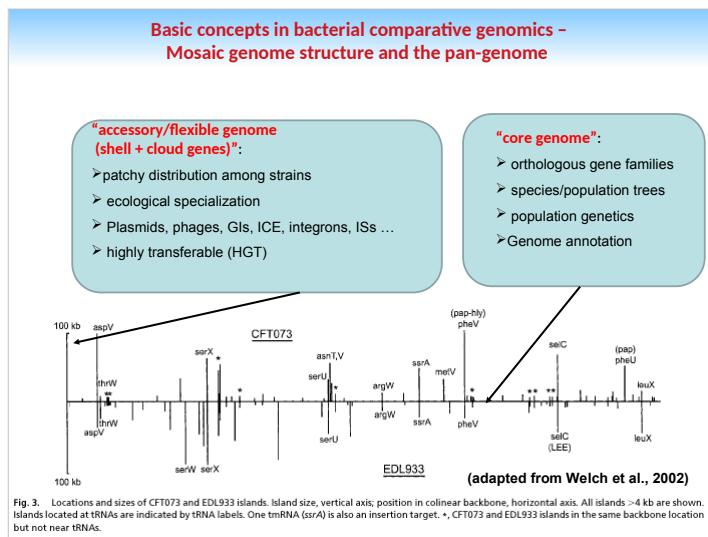
Contents

1. Basic concepts in bacterial comparative genomics – defining the pan-genome
 - genome mosaicism and the pan-genome
 - Structure of the pan-genome: the core and accessory/flexible (shell and cloud) genomes
 - homologous recombination and the genetic structure of bacteria populations
 - MGEs (ISs, GIs, Integrons, ICEs, prophages and plasmids) and the accessory genome
2. Computational challenges and approaches in pan-genome research
 - computational approaches to identify homologs: BDBH, COGs, OMCL
3. Popular software implementations to compute and analyze bacterial pan-genomes
 - 3.1 Computationally-intensive methods for accurate homology inference in diverged species
 - get_homologs
 - 3.2 Heuristic clustering methods for rapid species-level pan-genome analysis
 - Roary
4. Phylogenomic approaches for pan-genome analyses
 - 4.1 Computationally-intensive approaches:
 - get_phylomarkers maximum-likelihood core-genome phylogenies with optimal markers
 - get_phylomarkers and ML or parsimony pan-genome phylogenies
 - 4.2 Heuristic approaches to resolve population-level phylogenies
 - clustering pan-genome presence-absence distance matrix



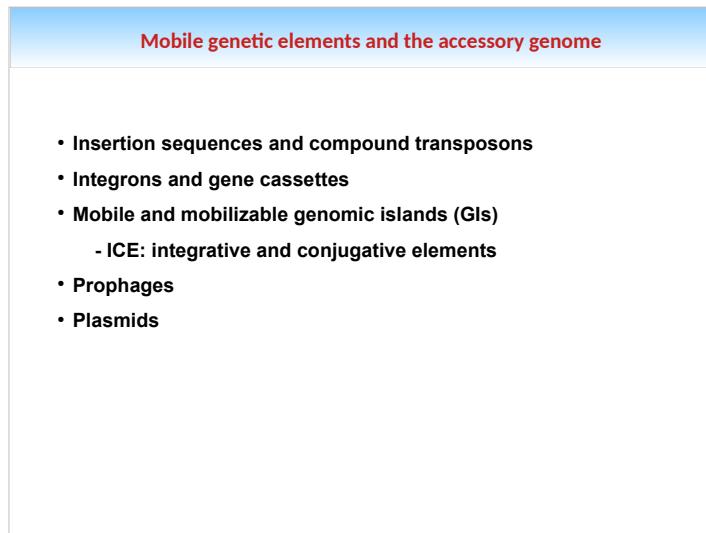
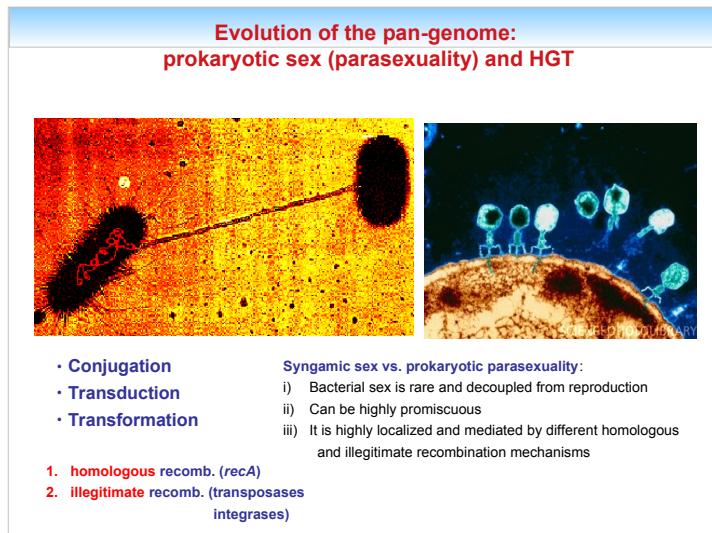
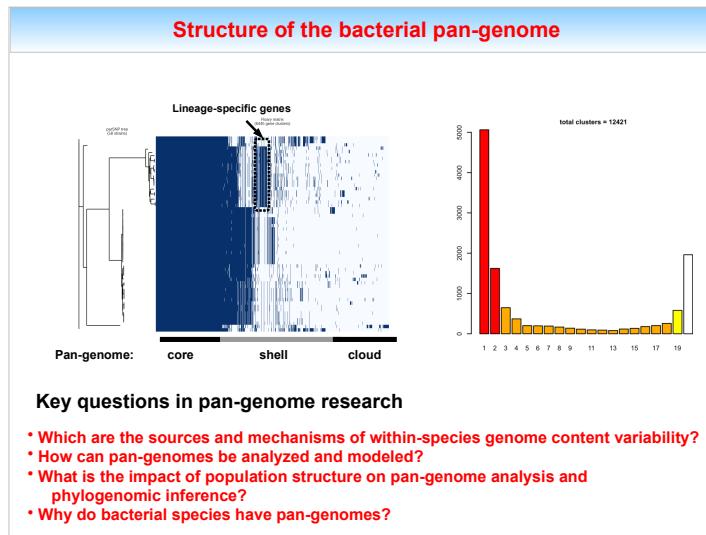
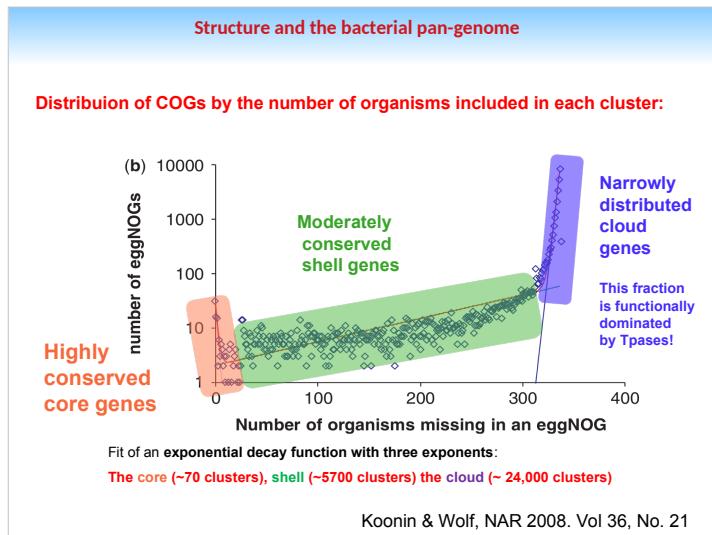
Introducción a la pan-genómica microbiana

IIIntroducción a la filoinformática – pan-genómica y filogenómica microbiana,
TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>



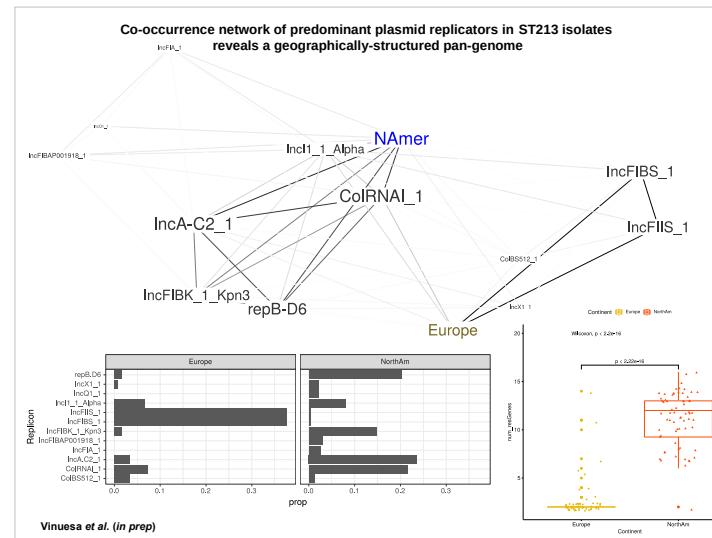
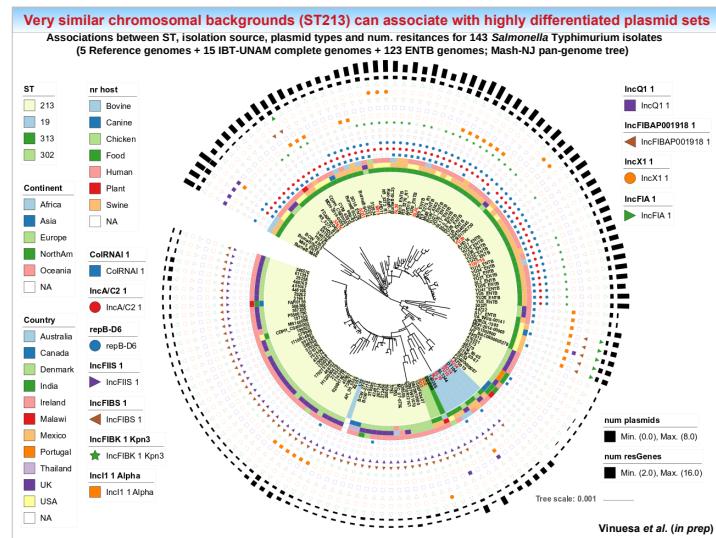
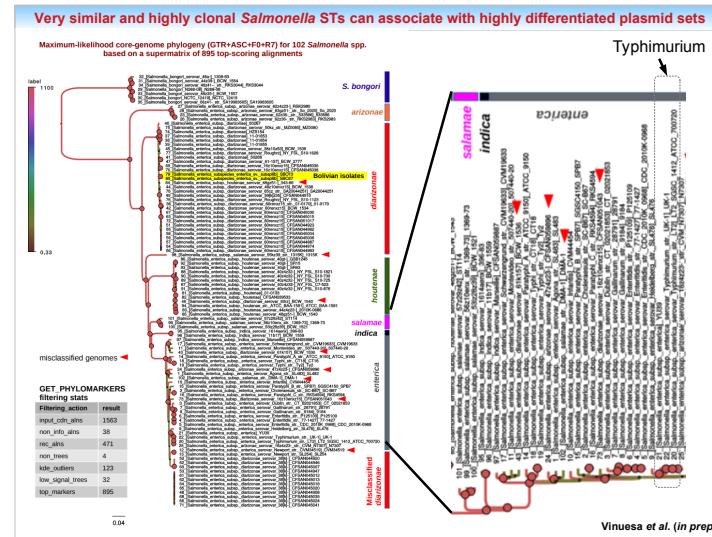
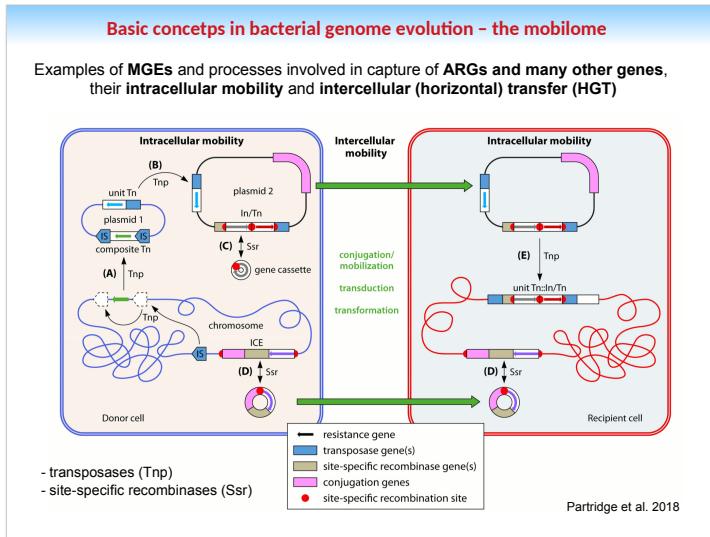
Introducción a la pan-genómica microbiana

IIIntroducción a la filoinformática – pan-genómica y filogenómica microbiana,
TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>



Introducción a la pan-genómica microbiana

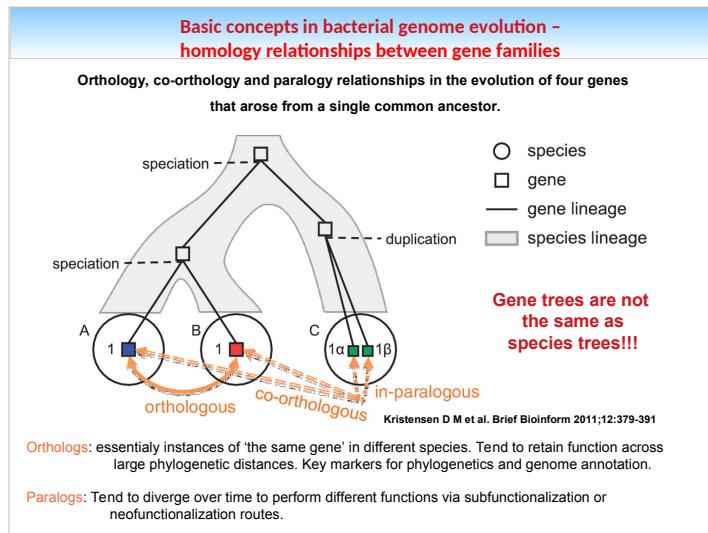
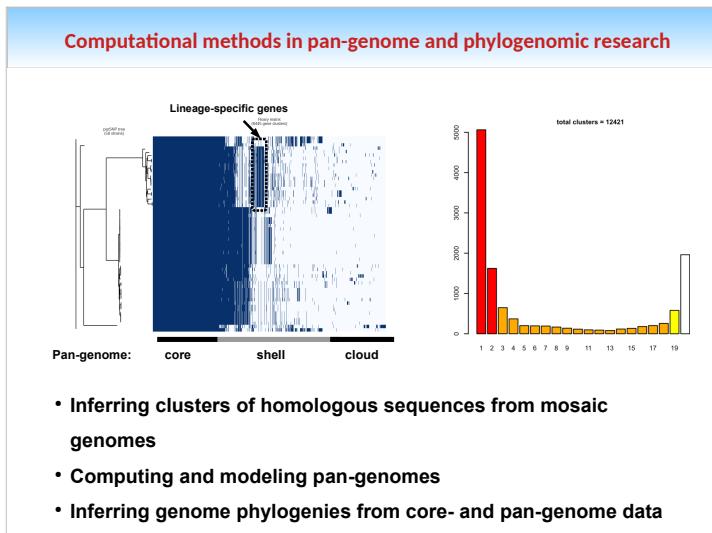
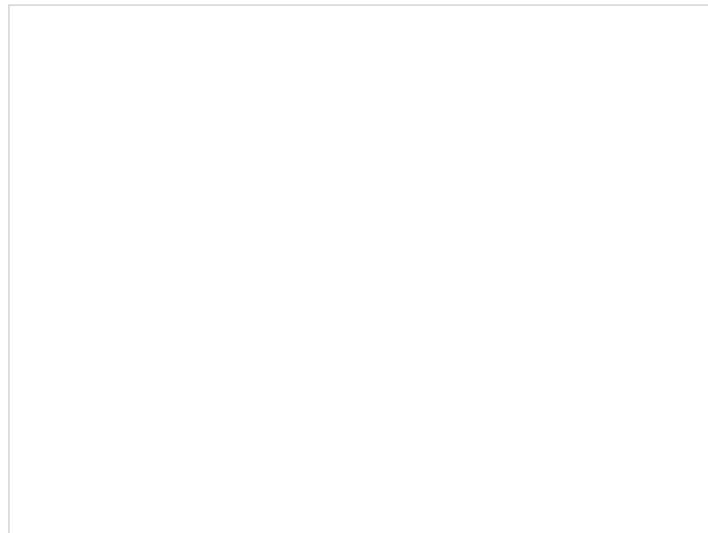
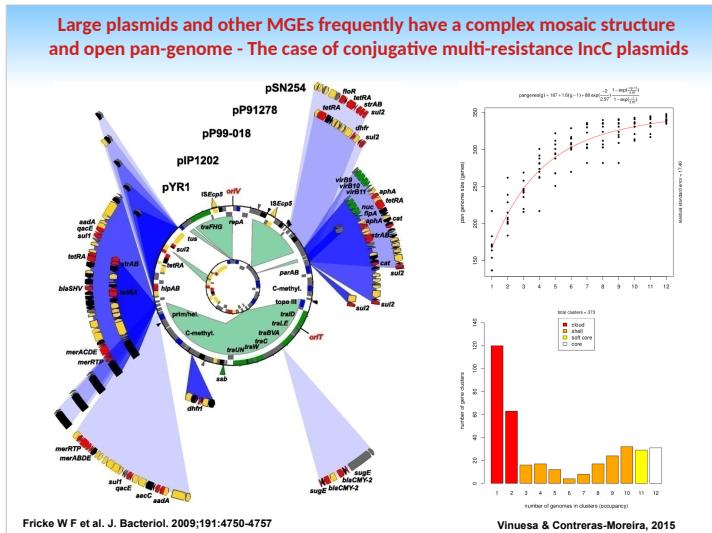
IIIntroducción a la filoinformática – pan-genómica y filogenómica microbiana,
TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>



© Pablo Vinuesa 2022. @pvinmex; vinuesa{at}ccg{dot}unam{dot}mx
<https://www.ccg.unam.mx/~vinuesa/>

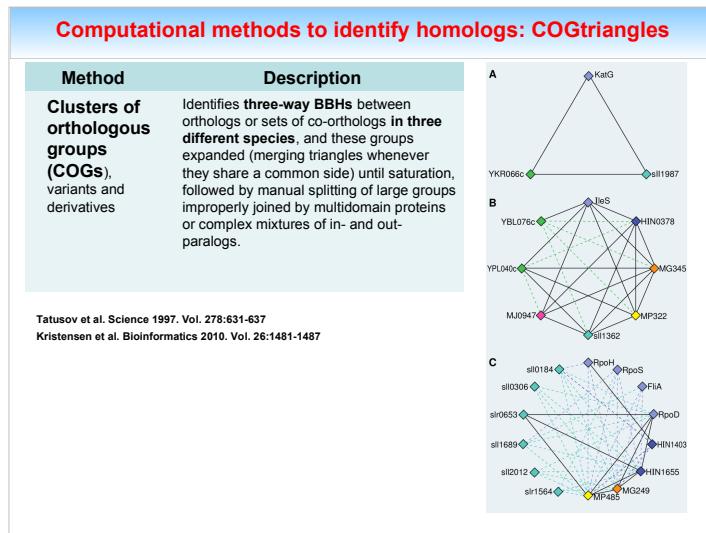
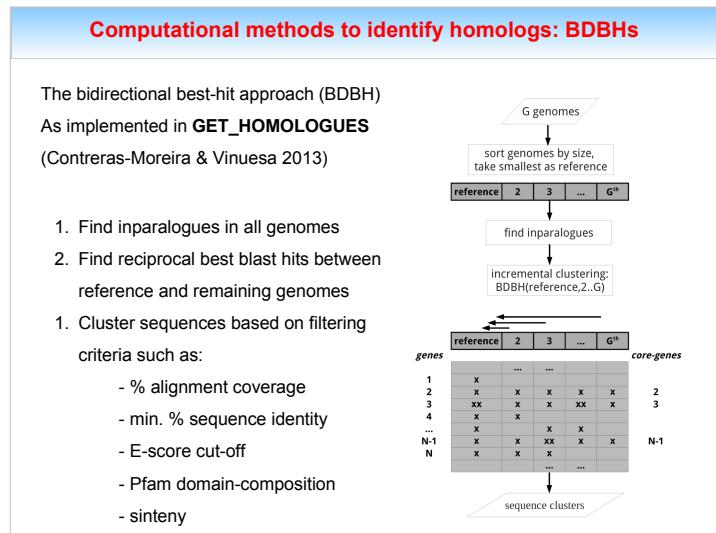
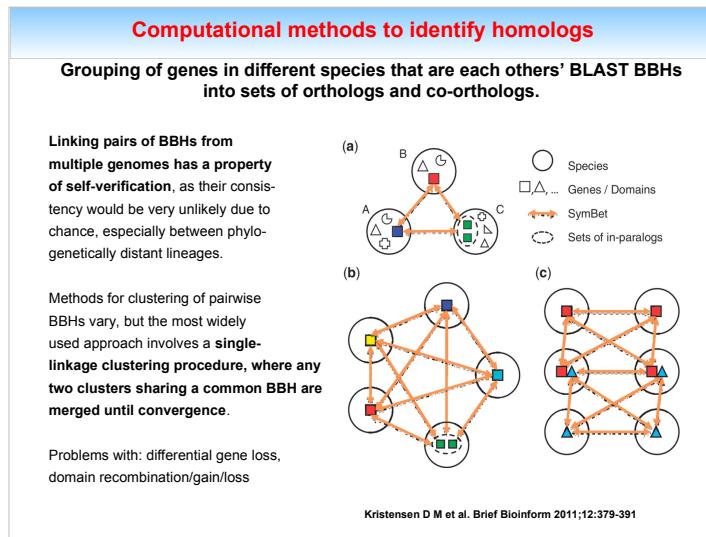
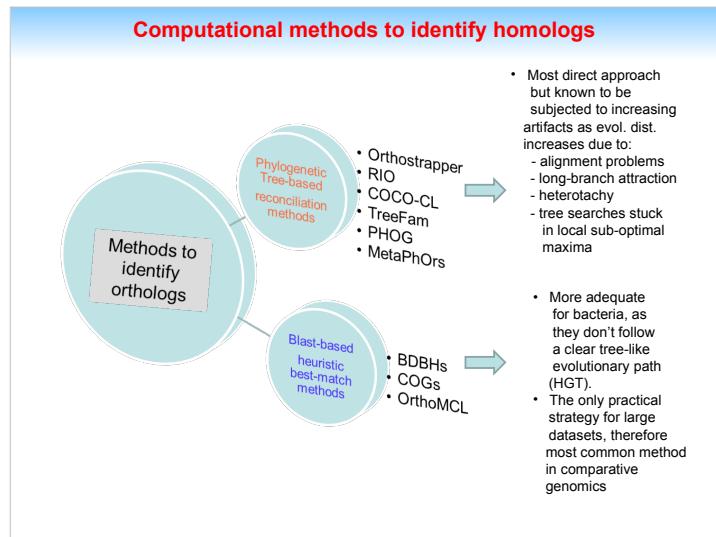
Introducción a la pan-genómica microbiana

IIIntroducción a la filoinformática – pan-genómica y filogenómica microbiana,
TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>



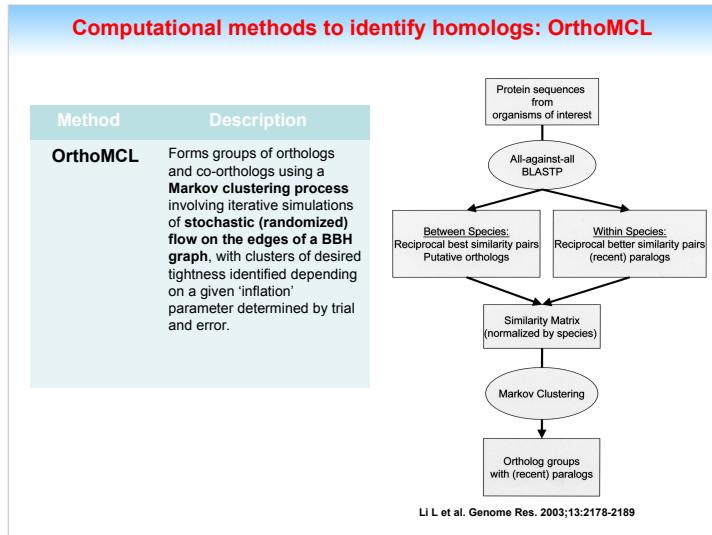
Introducción a la pan-genómica microbiana

Introducción a la filoinformática – pan-genómica y filogenómica microbiana,
TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>



Introducción a la pan-genómica microbiana

Introducción a la filoinformática – pan-genómica y filogenómica microbiana, TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>



Open-source tools for phylogenomics and microbial pan-genomics GET HOMOLOGUES & GET PHYLOMARKERS

Pablo Vinuesa
 I'm a microbiologist interested in genomics, ecology and evolution. I develop open-source code for these topics in collaboration with @eead-csic-complbio

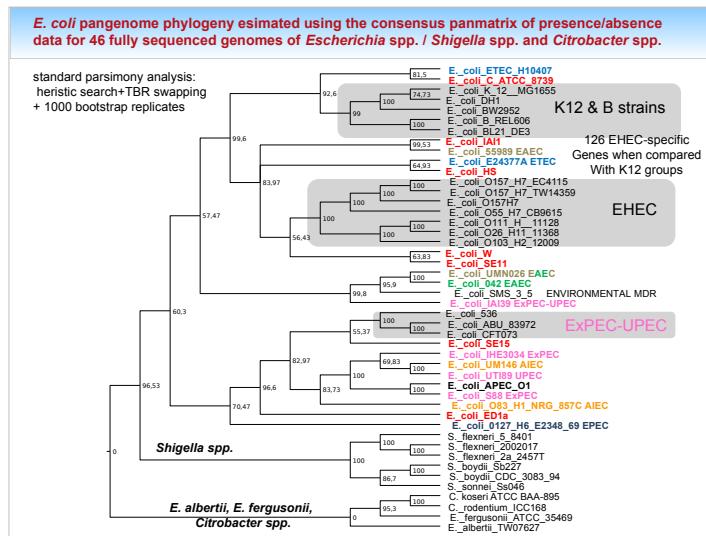
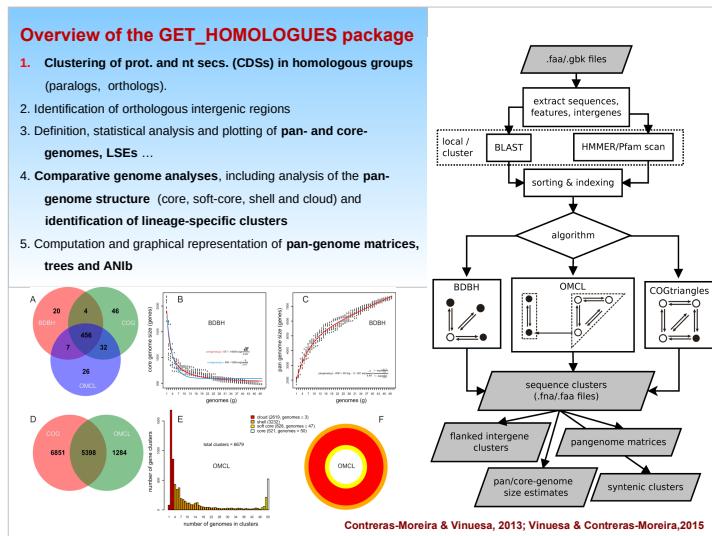
388 contributions in the last year

Contribution activity

Achievements

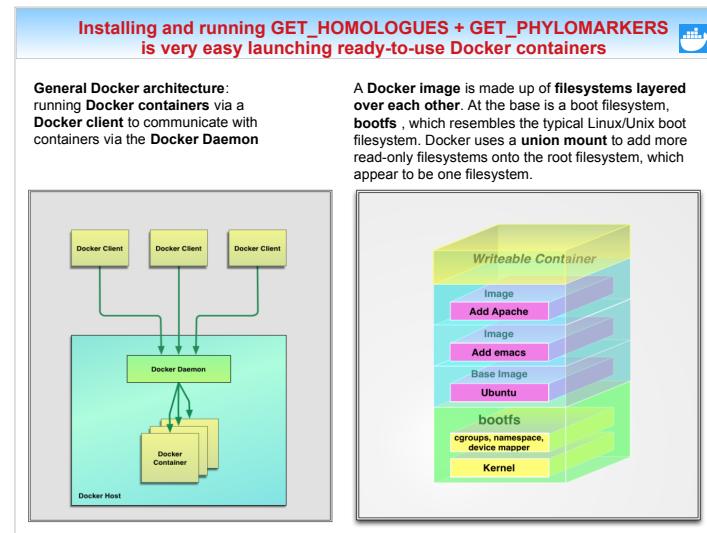
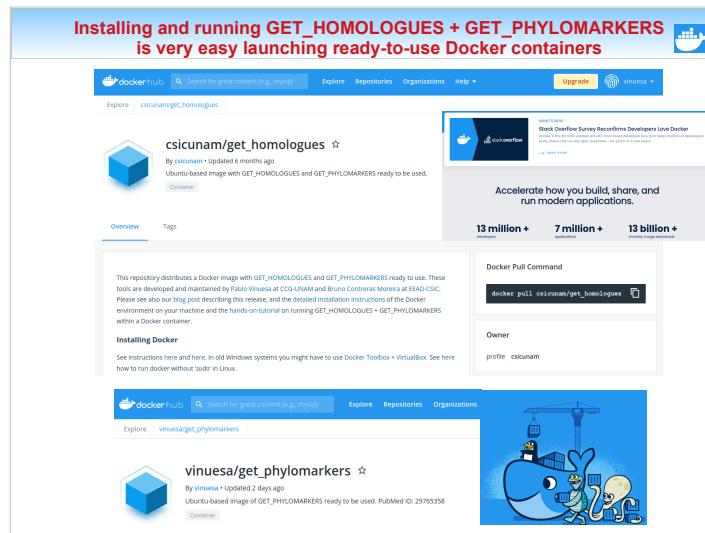
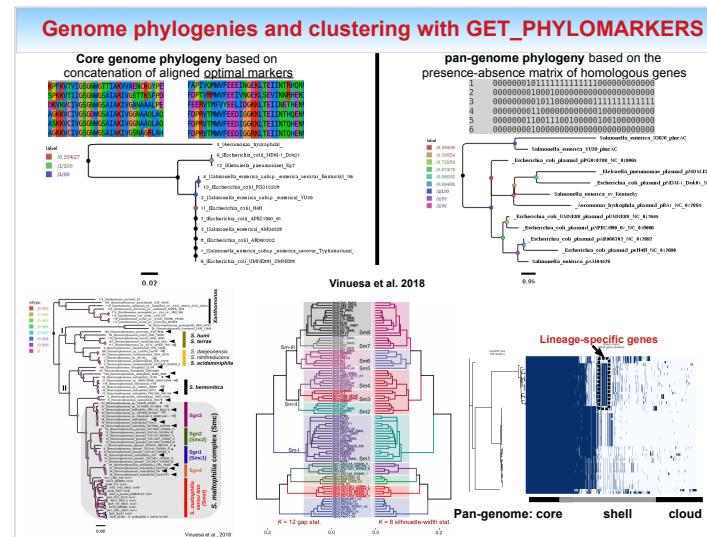
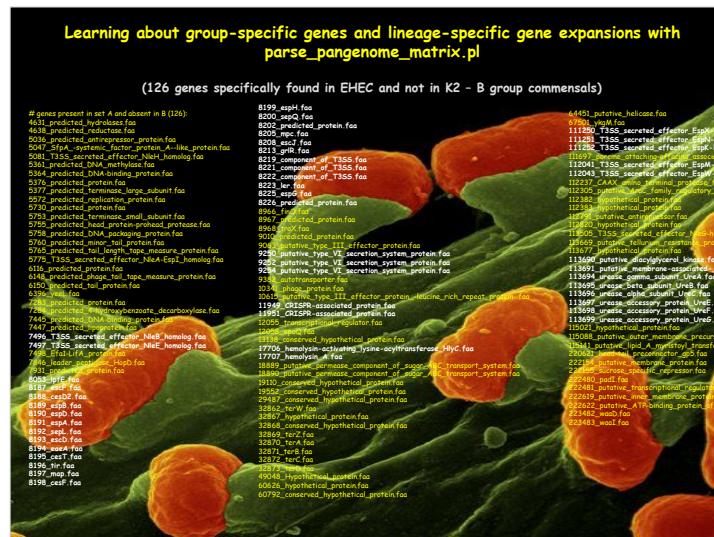
Contributions

<https://github.com/vinuesa>



Introducción a la pan-genómica microbiana

Introducción a la filoinformática – pan-genómica y filogenómica microbiana, TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>



Introducción a la pan-genómica microbiana

Introducción a la filoinformática – pan-genómica y filogenómica microbiana, TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>

Basic Docker commands and use of containers

```
## To avoid permission errors (and the use of sudo), add your user to the docker group
# https://docs.docker.com/install/linux/linux-postinstall/
sudo groupadd docker
sudo usermod -aG docker $USER

# 1. Get general docker info and print help
$ docker info
$ docker --help
$ docker run --help

# 2. List available docker images on your system
$ docker image ls

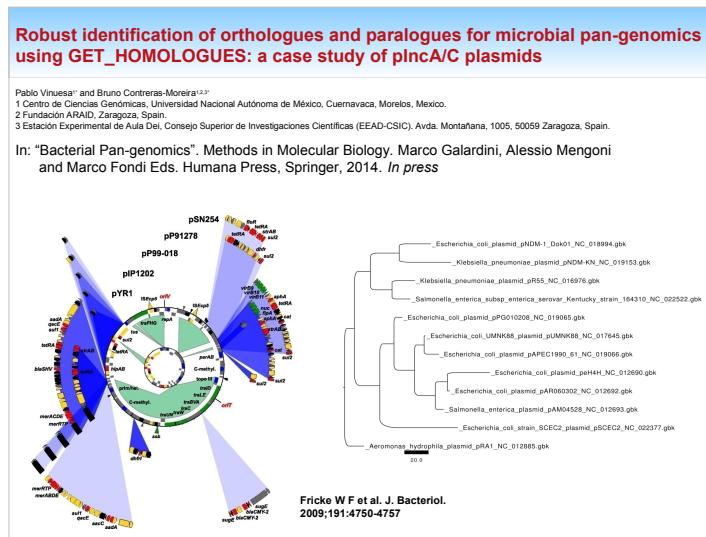
# 3. List running containers
$ docker container ls

# 4. Stop a container
$ docker container stop CONTAINER-ID

# 5. Pull a Docker image from the registry
$ docker pull csicunam/get_homologues:latest

# 6. Launch an image, using a mount-bind of a user directory on the Docker container
$ docker run --rm -it -v $HOME/get_homphy:/home/you/get_homphy \
  csicunam/get_homologues:latest /bin/bash

# The last command uses options --rm to remove the container after exiting and sets an
-i interactive session calling a pseudo tty (-t), mounting a host directory on the
container (-v ...), accessible for wr from both, and launching a bash shell (/bin/bash)
```



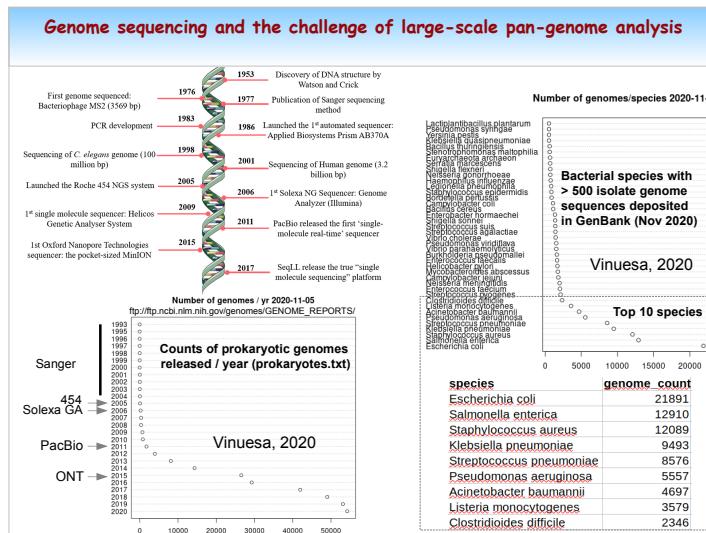
The GET_HOMOLOGUES + GET_PHYLOMARKERS tutorials:

- https://github.com/vinuesa/get_phylomarkers/

Analyses to be performed in an upcoming practical session:

A pangenomic analysis of *plncA/C* plasmids using GET_HOMOLOGUES

- Defining a robust core- and pan-genome of *plncA/C* plasmids
- Exploring the gene space of *plncA/C* plasmids: core, shell, cloud
- Pan-genome trees vs. core genome trees (supermatrices)
- Identifying lineage-specific genes in NDM-1 producing plasmids



Introducción a la pan-genómica microbiana

Introducción a la filoinformática – pan-genómica y filogenómica microbiana, TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>

Genome sequencing and the challenge of large-scale pan-genome analysis

Roary: rapid large-scale prokaryote pan genome analysis Bioinformatics, 31(22), 2015, 3691–3693

Andrew J. Page^{1,*}, Carla A. Cummins¹, Martin Hunt¹, Vanessa K. Wong^{1,2}, Sandra Reuter², Matthew T.G. Holden³, Maria Fookes⁴, Daniel Falush⁴, Jacqueline A. Keane¹ and Julian Parkhill¹

- Standard pan-genome software approaches

- The construction of a pan genome is **NP-hard** (Nguyen et al. 2014) with additional difficulties from real data due to contamination, fragmented assemblies and poor annotation. Therefore, any approach requires **heuristics** to produce a pan genome
- all-against-all** comparison using **BLAST**, with the **running time growing ~quadratically** with the size of input data and are computationally infeasible with large datasets.
- They also have **quadratic memory requirements**, quickly exceeding the RAM available in high performance servers for large datasets.

- Roary's heuristic approach

- Roary address the computational issues by performing a **rapid clustering of highly similar sequences with CD-HIT** which can reduce the running time of BLAST substantially, and carefully manages RAM usage so that it increases linearly, both of which make it possible to analyze datasets with thousands of samples using commonly available computing hardware.

Genome sequencing and the challenge of large-scale pan-genome analysis

- Roary's pipeline (Page et al. (2015) Bioinformatics, 31(22), 2015, 3691–3693)

- Prokka-annotated genomes (GFF3 + FASTA)
- Extract all the coding sequences, convert them into protein, and **cluster with CD-HIT**
- An **all-against-all** comparison is performed with **BLASTP** on the clustered sequences with a user defined percentage sequence identity (default 95%)
- Sequences are then clustered with **MCL** (Enright et al., 2002)
- Finally, the pre-clustering results from CD-HIT are merged together with the results of MCL
- Using conserved gene neighborhood information, homologous groups containing paralogs are split into groups of true orthologs
- A graph is constructed of the relationships of the clusters based on the order of occurrence in the input sequences, allowing for the clusters to be ordered and thus providing context for each gene.
- Isolates are clustered based on gene presence in the accessory genome, with the contribution of isolates to the graph weighted by cluster size.
- A suite of command line tools is provided to interrogate the dataset providing union, intersection and complement.

Genome sequencing and the challenge of large-scale pan-genome analysis

- Roary's pipeline benchmark (Page et al. (2015) Bioinformatics, 31(22), 2015, 3691–3693)

- 1 processor (AMD Opteron 6272), 60 GB of RAM.
- Constructed a simulated dataset based on *Salmonella Typhi* (*S.typhi*) CT18, to accurately assess the quality of the clustering.
- Created 12 genomes with 994 identical core genes and 23 accessory genes in varying combinations
- The running time and RAM of PGAP and PanOCT increases substantially with dataset size; LS-BSR, over clusters in 2% of cases.

Table 1. Accuracy of each pan genome application on a dataset of simulated data

	Core genes	Total genes	Incorrect split	Incorrect merge
Expected	994	1017	0	0
PGAP	991	1012	4	4
PanOCT	993	1015	1	1
LS-BSR	974	994	0	23
Roary	994	1017	0	0

Table 2. Comparison of pan genome applications using real *S.typhi* data (ERPP01718)

Samples	Software	Core ^a	Total	RAM (mb)	Wall time (s)
8	PGAP	4545	4929	669	41397
	PanOCT	4544	4936	663	1457
	LS-BSR	4474	4816	270	2385
	Roary	4459	4871	156	44
24	PGAP	—	—	—	—
	PanOCT	4522	4991	5313	96093
	LS-BSR	4451	4843	554	7807
	Roary	4445	4941	444	382
100	PGAP	—	—	—	—
	PanOCT	—	—	—	—
	LS-BSR	4272	7265	17413	345019
	Roary	4016	9201	13752	15465

^a Core is defined as the genes present in all samples, which allows for some assembly merges in very large datasets. Where there are no merges, the applications failed to complete within 5 days and used more than 40 GB of RAM. The first column is the number of unique *S.typhi* genomes in the input set with a mean of 54 contours over all 1000 assemblies.

Fig. 1. Effect of dataset size on the wall time of multiple applications. Only analysis that completed within 2 days and 60 GB of RAM is shown

Taking population-structure into account subdivides the pan-genome into 13 categories

MICROBIAL GENOMICS RESEARCH ARTICLE Horesh et al., Microbial Genomics 2021:7:000670 DOI 10.1093/mg/mab00670 MICROBIOLOGY DATA ACCESS

Different evolutionary trends form the twilight zone of the bacterial pan-genome

Gal Horesh¹, Alyce Taylor-Brown¹, Stephanie McGimpsey², Florent Lassalle¹, Jukka Corander^{1,2}, Eva Heinig^{2,*} and Nicholas R. Thomson^{1,3*}

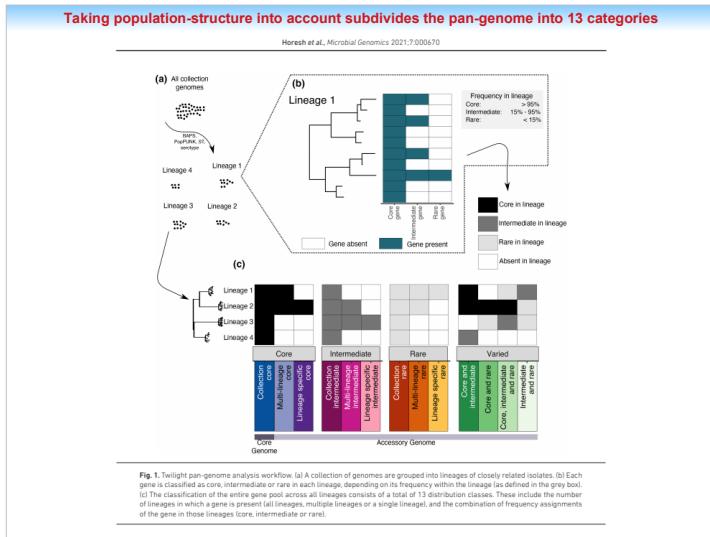
Abstract
The pan-genome is defined as the combined set of all genes in the gene pool of a species. Pan-genome analyses have been very useful in helping to understand different evolutionary dynamics of bacterial species: an open pan-genome often indicates a free-living lifestyle with metabolic versatility, while closed pan-genomes are linked to host-restricted, ecologically specialized bacteria. A detailed understanding of the species pan-genome has also been instrumental in tracking the phylogenetics of emerging drug resistance mechanisms. However, these analyses have focused on the core genes, ignoring the accessory genes. Here we present the application of a population structure-aware approach for classifying genes in a pan-genome based on within-species distribution. We demonstrate our approach on a collection of 7500 *Escherichia coli* genomes, one of the most-studied bacterial species and used as a model for an open pan-genome. We reveal clearly distinct groups of genes, clustered by different underlying evolutionary dynamics, and provide a more biologically informed and accurate description of the species' pan-genome.

The pan-genome has a heterogeneous distribution across the species' members. Traditionally, this gene pool has been divided into core genes, present across the majority of genomes, and accessory genes, whose presence vary across the dataset. These traditional methods do not reflect the true complexity of gene dynamics across a diverse species, nor do they account for population structure or the inherent biases in sampling and sequencing.

Horesh et al. propose a novel framework that further divides the core and accessory categories into 13 subcategories to better account for differences between lineages at a finer scale than traditional methods.

Introducción a la pan-genómica microbiana

Introducción a la filoinformática – pan-genómica y filogenómica microbiana, TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>



Why do prokaryotes have pan-genomes?

Why prokaryotes have pangenomes

James O. McInerney^{1*}, Alan McNally² and Mary J. O'Connell³

The existence of large amounts of within-species genome content variability is puzzling. Population genetics tells us that fitness effects on traits—such as differential survival or reproduction—coupled with the long-term effective population size of the species determines the likelihood of a new variant being removed, spreading to fixation or remaining polymorphic. Consequently, we expect that selection and drift will reduce genetic variation, which makes large amounts of gene content variation in some species so puzzling. Here, we amalgamate population genetic theory with models of horizontal gene transfer and assert that pangenomes most easily arise in organisms with large long-term effective population sizes, as a consequence of acquiring advantageous genes, and that the focal species has the ability to migrate to new niches. Therefore, we suggest that pangenomes are the result of adaptive, not neutral, evolution.

Nat. Microbiol. 28 MARCH 2017 | VOLUME: 2 | 17040

- There is a controversy around the question: are pangenomes adaptive?
 - HGT-mediated introgression of advantageous genes could lead to a **selective sweep**, reducing diversity
 - This would happen particularly in populations with small Ne (num. of individuals that contribute offspring to the next generation)
- An adaptive model for HGT when coupled with migration
 - A new compartment model by Niehus et al. that explicitly models HGT and migration has shown the plausibility of selection on HGT genes driving population differentiation. Using a mathematical approach, the authors showed that in the case of a selectively advantageous HGT event, diversity is removed from the species when there is no migration into or out of the compartment or niche occupied by the focal prokaryotic community.
- By contrast, a model that includes migration to and from the niche, combined with HGT of a selectively advantageous gene, can theoretically result in a situation where diversity is not necessarily reduced.
- There is a growing body of evidence that accessory genes might provide significant benefit. Karcagi et al. analysed a range of *E. coli* genomes at different levels of gene deletion, specifically genes that had been recently acquired by HGT. They found that HGT genes conferred significant benefits in terms of substrate utilization, efficiency of resource usage to build new cells, and tolerance to stress. Loss of HGT genes tended to affect fitness in several measurable ways, including the induction of a general stress response, inability to grow at all in some environments ...

Why do prokaryotes have pan-genomes?

Why prokaryotes have pangenomes

James O. McInerney^{1*}, Alan McNally² and Mary J. O'Connell³

The existence of large amounts of within-species genome content variability is puzzling. Population genetics tells us that fitness effects on traits—such as differential survival or reproduction—coupled with the long-term effective population size of the species determines the likelihood of a new variant being removed, spreading to fixation or remaining polymorphic. Consequently, we expect that selection and drift will reduce genetic variation, which makes large amounts of gene content variation in some species so puzzling. Here, we amalgamate population genetic theory with models of horizontal gene transfer and assert that pangenomes most easily arise in organisms with large long-term effective population sizes, as a consequence of acquiring advantageous genes, and that the focal species has the ability to migrate to new niches. Therefore, we suggest that pangenomes are the result of adaptive, not neutral, evolution.

Nat. Microbiol. 28 MARCH 2017 | VOLUME: 2 | 17040

Conclusion

- The authors infer that effective population size and the ability to migrate to new niches are the most influential factors in determining pan-genome size.
- Given the huge size of pan-genome estimates implies that the number of ecological niches on the planet must be enormous, which is in line with recent analysis of genomic diversity that suggests there may be 10^{12} (one trillion) microbial species on Earth.
- That the majority of genes in the biosphere are not strongly attached to any group of organisms has been a surprise of the genomics era, and consequently this 'public goods' hypothesis needed explanation

Bacterial evolution, mobile genetic elements and the pan-genome

References

- * Homology and genome evolution
Fitch WM. Homology: a personal view on some of the problems. *Trends Genet.* 2000 May;16(5):227-31
- Koonin EV. Orthologs, paralogs, and evolutionary genomics. *Annu Rev Genet.* 2005;39:309-38
- Koonin EV, Wolf YI. Genomics of bacteria and archaea: the emerging dynamic view of the prokaryotic world. *Nucleic Acids Res.* 2008 Dec;36(21):6688-719.
- Kristensen DM, Wolf YI, Mushegian AR, Koonin EV. Computational methods for Gene Orthology inference. *Brief Bioinform.* 2011 Sep;12(5):379-91. <=
- * Bacterial mobile genetic elements
Aziz RK, Breitbart M, Edwards RA. Transposases are the most abundant, most ubiquitous genes in nature. *Nucleic Acids Res.* 2010 Jul;38(13):4207-17
- Bellanger X, Payot S, Leblond-Bouget N, Guédon G. Conjugative and mobilizable genomic islands in bacteria: evolution and diversity. *FEMS Microbiol Rev.* 2014 Jul;38(4):720-80
- Cambray G, Gueout AM, Mazel D. Integrins. *Annu Rev Genet.* 2010;44:141-66
- Campbell A. The future of bacteriophage biology. *Nat Rev Genet.* 2003 Jun;4(6):471-7.
- Frost LS, Leplae R, Summers AO, Toussaint A. Mobile genetic elements: the agents of open source evolution. *Nat Rev Microbiol.* 2005 Sep;3(9):722-32
- Galán JE, Waksman G. Protein-Injection Machines in Bacteria. *Cell.* 2018 Mar 8;172(6):1306-1318
- Gillings MR. Integrins: past, present, and future. *Microbiol Mol Biol Rev.* 2014 Jun;78(2):57-77
- Johnson CM, Grossman AD. Integrative and Conjugative Elements (ICEs): What They Do and How They Work. *Annu Rev Genet.* 2015;49:577-601
- Juhás M, van der Meer JR, Gaillard M, Harding RM, Hood DW, Crook DW. Genomic islands: tools of bacterial horizontal gene transfer and evolution. *FEMS Microbiol Rev.* 2009
- Partridge SR, Kwong SM, Firth JN, Jensen SO. Mobile Genetic Elements Associated with Antimicrobial Resistance. *Clin Microbiol Rev.* 2018 Aug 1;31(4):<=
- Pilla G, Tang CM. Going around in circles: virulence plasmids in enteric pathogens. *Nat Rev Microbiol.* 2018 Aug;16(8):484-495 <=

Introducción a la pan-genómica microbiana

IIIntroducción a la filoinformática – pan-genómica y filogenómica microbiana,
TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>

Bacterial evolution, mobile genetic elements and the pan-genome

References

* Bacterial mobile genetic elements (continuation)

- Siguier P, Courteyre E, Varani A, Ton-Hoang B, Chandler M. Everyman's Guide to Bacterial Insertion Sequences. *Microbiol Spectr*. 2015 Apr;3(2):MDNA3-0030-2014.
- Welch RA, Burland V, Plunkett G 3rd, Redford P, Roesch P, Rasko D, Buckles EL, Liou SR, Boutin A, Hackett J, Stroud D, Mayhew GF, Rose DJ, Zhou S, Schwartz DC, Perna NT, Mobley HL, Donnenberg MS, Blattner FR. Extensive mosaic structure revealed by the complete genome sequence of uropathogenic Escherichia coli. *Proc Natl Acad Sci U S A*. 2002 Dec 24;99(26):17020-4.
- Wozniak RA, Waldor MK. Integrative and conjugative elements: mosaic mobile genetic elements enabling dynamic lateral gene flow. *Nat Rev Microbiol*. 2010 Aug;8(8):552-63.

Bacterial evolution, mobile genetic elements and the pan-genome

References

* Mosaic genome structure, population structure and the bacterial pan-genome: reviews and papers

- Golicz, A. A., Bayer, P. E., Bhalla, P. L., Batley, J. & Edwards, D. Pangenomics Comes of Age: From Bacteria to Plant and Animal Applications. *Trends in Genetics* vol. 36 132–145 (2020). <=
- Horesh et al. 2021. Different evolutionary trends form the twilight zone of the bacterial pan-genome. *Microb Genom*. 2021 Sep;7(9). <=
- McInerney JO, McNally A, O'Connell MJ. Why prokaryotes have pangenomes. *Nat Microbiol*. 2017 Mar 28;2:17040 <=
- Medini D, Donati C, Tettelin H, Masignani V, Rappoli R. The microbial pan-genome. *Curr Opin Genet Dev*. 2005 Dec;15(6):589-94.
- Tettelin H, Masignani V, Cieslewicz MJ, Donati C, Medini D, et al. Genome analysis of multiple pathogenic isolates of Streptococcus agalactiae: implications for the microbial "pan-genome". *Proc Natl Acad Sci U S A*. 2005 Sep 27;102(39):13950-5. doi: 10.1073/pnas.0506758102. Epub 2005 Sep 19. Erratum in: Proc Natl Acad Sci U S A. 2005 Nov 8;102(45):16530. PMID: 16172379; PMCID: PMC1216834.
- Tettelin H, Medini D, editors. *The Panogenome: Diversity, Dynamics and Evolution of Genomes*. Cham (CH): Springer; 2020. PMID: 32633920.
- Welch RA, Burland V, Plunkett G 3rd, Redford P, Roesch P, Rasko D, Buckles EL, Liou SR, Boutin A, Hackett J, Stroud D, Mayhew GF, Rose DJ, Zhou S, Schwartz DC, Perna NT, Mobley HL, Donnenberg MS, Blattner FR. Extensive mosaic structure revealed by the complete genome sequence of uropathogenic Escherichia coli. *Proc Natl Acad Sci U S A*. 2002 Dec 24;99(26):17020-4.

* Tools and approaches for computing pan-genomes and (pan-)genome phylogenies

- Contreras-Moreira B, Vinuesa P. GET_HOMOLOGUES, a versatile software package for scalable and robust microbial pangenome analysis. *Appl Environ Microbiol*. 2013 Dec;79(24):7696-701.
- Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MT, Fookes M, Falush D, Keane JA, Parkhill J, Roarty: rapid large-scale prokaryote genome analysis. *Bioinformatics*. 2015 Nov 15;31(22):3691-3.
- Vinuesa P, Contreras-Moreira B. Robust identification of orthologues and paralogues for microbial pan-genomics using GET_HOMOLOGUES: a case study of plncA/C plasmids. *Methods Mol Biol*. 2015;1231:203-32.
- Vinuesa P, Ochoa-Sánchez LE, Contreras-Moreira B. GET_PHYLOMARKERS, a Software Package to Select Optimal Orthologous Clusters for Phylogenomics and Inferring Pan-Genome Phylogenies. Used for a Critical Geno-Taxonomic Revision of the Genus Stenotrophomonas. *Front Microbiol*. 2018 May 1:9771.